# EXHIBIT F

(Part 2)

1    events, recorded some information which would lead

2    him to conclude that the robots.txt exclusion

3    pages between 4500 and 5,000 bytes?  Do you know

4    that?

5         A.   I don't know whether that's the case or

6    not, but that's not what I was saying before.  I

7    was talking about whether there were other files

8    of the same size.

9         Q.   Yes.  That's not my question.  My

10   question is a little bit more precise.  I am

11   asking you to consider whether the bases for

12   Mr. Lenky's analysis -- and I'll go through them

13   in more exhaustive detail if you like -- number

14   one, that robots.txt exclusion pages are between

15   4500 and 5,000 bytes, number one, and that by

16   filtering out the HTTP access logs for successful

17   data transmissions between Internet Archive and

18   the defendants where the transmission size was in

19   range, that same range, and then exclude image

20   files and the robots.txt file, a subset of log

21   file entries representing potential robots.txt

22   exclusion page views, etc. -- whether, given that

23   this information was based upon Mr. Mohr and his

24   knowledge, whether that would change your concern

25   about the conclusion that was reached.

1          MR. LEWIS:  Objection to the form.

2      A.    No.   What I'm saying is that the numbers

3  that Mr. Lenky concludes here do not follow from

4  the assumptions and assert -- the facts he claims

5  here in this paragraph.  There are other things

6  one would have to assume to reach -- to reach the

7  numbers that he reaches, and he does not state

8  those other assumptions, let alone giving any

9  reason for why they would be true.

10     Q.    Well, what, in your view, are the

11 fallacies with his reasoning in this area?

12     A.    The fallacy is assuming -- one fallacy

13 is assuming that there could not have been any

14 file between 4,500 and 5,000 bytes in size other

15 than the ones that he specifically excluded.

16     Q.    What do you mean more precisely than

17 that?

18     A.    Well, he gets to his final numbers by

19 subtraction.  He first identifies all of the

20 accesses that were in that size range, and then he

21 subtracts off some files that were of that size

22 range that are not -- that are not exclusion

23 pages.  In order to get to the accurate number, he

24 would have to exclude all of the pages of that

25 size that were not exclusion files.  He subtracts

1    off some of them, but he --

2        Q.    But you're not -- I don't --

3        A.    -- but he doesn't say that he's

4    subtracting off all of them.  He doesn't give any

5    reason for believing that he subtracted off all of

6    them.  There are other categories of pages that he

7    has not excluded here that would be included in

8    this final number.

9        Q.    Anything else?

10       A.    Well, that's certainly the primary --

11   the primary reason to question his numbers here.

12       Q.    Well, if there are any secondary

13   reasons, I want to know about them.  Is there

14   anything else?

15       A.    Nothing else comes to mind.

16       Q.    Looking at the next paragraph -- if you

17   want to take a moment to read it, please do.  Have

18   you read through that, the next paragraph?

19       A.    Let me take a minute.

20   Okay.

21       Q.    Do you have any dispute as to the

22   methodology that he used in that paragraph to

23   confirm the results and conclusions from the prior

24   paragraph?

25       A.    I understand the methodology he used.

1       Q.   Did you otherwise, through any other

2   source, come to a conclusion as to which

3   individuals may have been directly involved?

4       A.   No.

5       Q.   I think you testified earlier that you

6   did not speak to anyone at Harding Earley at any

7   time for any purpose.  Correct?

8       A.   That's correct.

9       Q.   Do you disagree with Mr. Lenky's

10  conclusion that the two distinct user-agent

11  strings present in the web logs indicate that

12  there were at least two separate machines who were

13  responsible for initiating their requests for the

14  content?

15      A.   I agree that different user-agent

16  strings would tend to indicate different

17  computers.  As to whether there were two different

18  strings, I would have to look to the logs.

19      Q.   Look at the logs.

20      A.   In a quick scan of the logs, I'm only

21  seeing one, but perhaps there's information you

22  can point me to.

23      Q.   I'll be happy to.  Can I have that

24  exhibit, please?

25      A.   Sure.

1      Q.    Referring you to the first page of the

2   exhibit, the first log entry and the third full

3   log entry from the bottom.

4      A.    Yes, I see those.

5      Q.    Those would be two distinct user-agent

6   strings?

7      A.    Yes, those are different user-agent

8   strings.

9      Q.    That would seem to suggest two different

10  machines were involved?

11     A.    Yes, that would tend to indicate that.

12     Q.    Look at page 43 of your report.

13     A.    You mean paragraph 43?

14     Q.    Yes, paragraph 43.  Sorry.  You mention

15  in this paragraph that they, the Harding Earley

16  employees, did not engage in "hacking."  How do

17  you define "hacking" in this context?

18     A.    Well, I'm using the term here to mean

19  activity that is -- activity that is devious and

20  out of the ordinary, essentially.

21     Q.    In your mind, hacking connotes an intent

22  to be devious?

23     A.    Well, the term "hacking" is used in many

24  different ways within computer science.  It's

25  sometimes used to refer only to someone exhibiting

1    following testimony:

2            "QUESTION:  Is it inconsistent with your

3    review of the logs that the following happened,

4    number one, that members of the Harding Earley law

5    firm repeatedly attempted to access an archived

6    web page of Healthcare Advocates within a short

7    span of time?")

8        A.    They did request archived pages from the

9    Internet Archive.  And depending on how short a

10   span of time and how many repetitions you're

11   talking about, I suppose you might be able to say

12   that they requested -- they made repeated requests

13   over a period of time.

14       Q.    Over a short period of time?

15            MR. LEWIS:  Objection.

16       Q.    When I say a "short period of time," I

17   mean a matter of minutes.  Did you observe that

18   through the logs?

19       A.    There are points in the log where one

20   sees more than one request in a period of a few

21   minutes.

22       Q.    For the same archived content.  Correct?

23       A.    For the same archived content, I'd need

24   to review the logs to answer that.

25       Q.    Please do.

1          Have you had a chance to review the logs,

2    Professor?

3          A.    Yes.

4          Q.    Can you answer the question now?

5          A.    I do see some instances of the same

6    content being requested close together in time.

7          Q.    Within a matter of minutes.  Correct?

8          A.    Yes.

9                MR. LEWIS:  Objection to form.

10         A.    Yes.

11         Q.    And do you also note that there's a

12   pattern throughout the logs whereby there are a

13   number of unsuccessful requests for particular

14   content followed by a successful request?

15               MR. LEWIS:  Objection to form.

16         A.    I don't know whether that's in the logs

17   or not, and I don't know that I can answer the

18   question at all by looking at printed-out logs.

19         Q.    What else would you need to review in

20   order to rely on that?

21         A.    I think I would need to use

22   pattern-matching tools and so on, to have an

23   electronic version of the logs and to operate on

24   that.

25         Q.    Did you use your pattern-matching tools

1    examples.  Generally trying to classify the

2    requests and to count different -- the different

3    kinds of events.

4         Q.   During the break when you were reviewing

5    the logs in order to answer the question that I

6    initially asked you, I believe you said that you

7    observed that there were instances within a matter

8    of minutes where there were repeated requests for

9    the same archived web content.  I don't want to

10   misstate your testimony, but is that accurate?

11        A.   More than one request for the same

12   content during a matter of a few minutes, yes.

13        Q.   Did you also observe that at the end of

14   that state of requests, there was a successful

15   access of that same requested web content?

16        A.   No.

17        Q.   What did you observe?

18        A.   I didn't try to look for the end of the

19   chain or the sequence.

20        Q.   Could you do that if you sat here today?

21        A.   It would take some time.  The logs in

22   this form are not easy to digest.

23        Q.   Do you recall at any point in time in

24   reviewing the logs seeing at least once that same

25   pattern, repeated attempts to access a specific

1    archived web page which were unsuccessful,

2    followed by a success in achieving access to that

3    web page within a short period of time, meaning a

4    span of minutes?

5          A.    No, I don't recall seeing that, that

6    pattern.

7          Q.    Did you specifically look for that

8    pattern?

9          A.    I don't recall looking for it.

10         Q.    Let me show you again what's been marked

11   Bonini-8, the robots exclusion page.  At the same

12   time I'll refer you to your report, paragraph 49

13   and 50.

14         Starting at paragraph 49, it's your opinion

15   that the law firm's accesses to the Wayback

16   Machine were not unauthorized, and, moreover, that

17   there was no intentional unauthorized access on

18   their part?

19         A.    Well, I disagree with Mr. Lenky's

20   apparent assertion that the law firm employees

21   knew that what they were doing was unauthorized.

22         Q.    Okay.  Well, then, I understand what

23   you're saying.  I think that's different from what

24   I'm saying.  I'm asking you two discreet

25   questions, and I'll ask them to you separately in

1          Q.    I'll rephrase the question if you can't

2      answer it.

3          A.    Well, it does appear that Healthcare

4      Advocates during at least part of the relevant

5      time had a robots.txt file in place that asked the

6      Internet Archive crawler not to visit their page.

7      As to what their intention was, I don't know.

8          Q.    Is it also clear to you that the robots

9      exclusion that Healthcare Advocates had in place

10     was effective in blocking third-party access to

11     the archived web content during that period of

12     time?

13              MR. LEWIS:  Objection to form.

14         A.    Was it effective?  Well, we know that

15     some accesses to the archived content were

16     successful.

17         Q.    But don't we also know that --

18              MR. LEWIS:  I don't think he's finished

19     his answer.  Have you finished your answer?

20              MR. CHRISTIE:  All right.  I'll let him

21     finish.

22         A.    I'll stop there.

23         Q.    Well, don't we also know that some

24     attempted accesses were rebuffed?

25         A.    Some attempted accesses were

1    unsuccessful, yes.

2         Q.    And that was based upon the robots

3    exclusion.  Correct?

4              MR. LEWIS:  Objection to the form.

5         A.    Well, that was the behavior of the

6    Wayback Machine when those requests were made.

7         Q.    So would you agree with me that at least

8    in part of the time during the period July 9,

9    2003, through July 14, 2003, it was a properly

10   configured and properly implemented robots.txt

11   exclusion on the Healthcare Advocates web server?

12             MR. LEWIS:  Objection to form.

13        A.    This was probably the case for at least

14   part of that period.

15        Q.    And is it possible that it could have

16   been that case for the whole period of time?

17        A.    That is possible.

18        Q.    Is it also accurate to say that, based

19   on your report, you dispute Mr. Lenky's conclusion

20   that any access -- that any unauthorized access,

21   to the extent that there is any, by the law firm

22   representatives was intentional?

23             MR. LEWIS:  Objection to the form.

24        A.    Well, Mr. Lenky doesn't explain his

25   reasoning, if I recall correctly, for why he says

1    file itself and the query exclusion page, which we

2    marked as Bonini-8.  Is that accurate?

3         A.    Yes, that's roughly what my report says.

4         Q.    What you don't talk about in your

5    report, number one, is the number of instances

6    that the Harding Earley representatives actually

7    physically and intentionally sought out and viewed

8    the robots.txt file, do you?

9              MR. LEWIS:  Objection to form.

10        A.    Well, I do say that he makes -- one of

11   the bases he appears to rely on is Harding Earley

12   employees seeing the robots.txt file.  And if

13   he -- and to the extent that he relies on how many

14   times that happened, that's something that I'm

15   discussing here.

16        Q.    Well, please look at page 7 of his

17   report.  There's a table that goes through to page

18   9 which sets forth specific instances where the

19   Harding Earley representatives are requesting the

20   robots.txt file.  Do you see that table?

21        A.    Yes.  He points to one user doing it

22   once and another user, he says, doing it seven

23   times.

24        Q.    Did you have occasion to, in the course

25   of your analysis, to compare the claims in this

1    table with the logs themselves?

2         A.    I don't recall whether I did or not.

3         Q.    Do you have any reason to dispute the

4    accuracy of the information in this table?

5         A.    Sitting here today, no.

6         Q.    Do you see that the -- that in the table

7    in Mr. Lenky's report, that the timing of the

8    requests for the robots.txt file are early in the

9    day on July 9?  When I say "early in the day," I

10   mean between 6:25 Pacific time and 8:09 Pacific

11   time.  Do you see that?

12        A.    Well, I see the times.  I'm not positive

13   about the time zone.  It says "minus 0700" on the

14   time zone.

15        Q.    Well, I thought you didn't dispute the

16   accuracy of the table, which includes --

17        A.    Well, there are times listed here.  It

18   says, for example, on the first access the time as

19   06:25:20.  Then it says "minus 0700," which I know

20   denotes a time zone, but I'm not sure which time

21   zone it denotes.

22        Q.    Okay.  Do this, please.  Take the logs,

23   ones for Internet Archive, which are marked

24   Mohr-2, and please look and see if you can find

25   the first referenced access in the logs at 6:25,

1     6:25:20.

2          A.    Yes, I see it.

3          Q.    Do you want to take a look to see if you

4     can find the other ones, too?

5          A.    Sure.  Are you asking me to confirm the

6     things on the right-hand side of the table or that

7     there accesses to robots.txt at the time listed on

8     the left?

9          Q.    While you're at it, I want you to

10    confirm everything.

11         A.    Okay, I've reviewed all of these.

12         Q.    Having reviewed in Mohr-2 the items in

13    the table on pages 7 through 9 of Mr. Lenky's

14    report, do you dispute any of the information

15    contained in the table as per your comparison with

16    the logs?

17         A.    There is, I think, one error in the

18    table --

19         Q.    What's that?

20         A.    -- which is in the second access for

21    user B.

22         Q.    Yes.

23         A.    I believe he has the time slightly

24    wrong.  It's off by a few seconds.  Also, I should

25    note that the language above the table that says

1    that "Each of these is an access to the robots.txt

2    file" is not correct.  The last entry in the table

3    is not a reference to the robots.txt file.  The

4    last entry in the table is not a request for the

5    robots.txt file.  It's a request for other

6    content.

7            Q.    Isn't that what the explanatory note --

8            A.    Yes, his explanatory note says that, but

9    the caption above the table says that -- this says

10    that the table is a list of requests for

11    robots.txt file, so that doesn't match.

12    Otherwise, these appear to be as described in the

13    table.

14            Q.    You agree with his interpretation of the

15    significance of the log entries in the right-hand

16    column.  Correct?

17            MR. LEWIS:  Objection to the form.

18            A.    His descriptions in the right-hand

19    column appear to be accurate.  As to what

20    interpretation he makes of that, that's perhaps a

21    different story.  Focusing on just he wrote in the

22    right-hand column, yes, that does appear to be

23    consistent with the logs.

24            Q.    So you would agree, then, that early in

25    the morning on July 9 of 2003 there were repeated

1      requests specifically for the robots.txt file

2      emanating from the Harding Earley law firm?

3                MR. LEWIS:  Objection to form.

4          A.   There were requests for the robots.txt

5      file from the Internet Archive.

6          Q.   And, in fact, there were requests not

7      only for the current version of the robots.txt

8      file.  Correct?

9                MR. LEWIS:  Same objection.

10         A.   Correct.

11         Q.   There were requests for versions from

12     the year 2000.  Correct?

13         A.   Yes, there appears to be the requests

14     for the version from the year 2000.

15         Q.   As well as a request for all archived

16     versions of the robots file.  Correct?

17               MR. LEWIS:  Same objection.

18         A.   Yes, there appears to be such a request.

19         Q.   So that would appear to indicate that

20     there was at least one user agent who was

21     conducting an investigation into the robots.txt

22     file?

23               MR. LEWIS:  Objection to the form.

24         A.   There was this one user agent who --

25     which might or might not have been a single person

1    or more than one person.  It appears to be one

2    computer from which these requests were made.

3         Q.    Does it also appear that the requests

4    were made for the robots text file from the query

5    exclusion page, Bonini-8?

6              MR. LEWIS:  Objection to form.

7         A.    That I can't be certain of.

8         Q.    Why can't you be certain of that?

9         A.    Well, the question is whether the

10   previous page that the user saw was an exclusion

11   page or not, and that is not immediately evident

12   from the log.

13        Q.    Okay.  Well, if you look at the text

14   string in the bottom left-hand corner of Bonini-8

15   and you compare it to the right-hand column of the

16   table in Mr. Lenky's report, would that alter the

17   answer you just gave me?

18        A.    No.

19        Q.    Why not?

20        A.    Because the requests that are made for

21   archived content would yield either an exclusion

22   page for either content or a list of content

23   that's available.  And you cannot immediately tell

24   from looking at the URL, which was provided in a

25   given instance.  So I can't tell, sitting here,

## IN THE UNITED STATES DISTRICT COURT
## FOR THE EASTERN DISTRICT OF PENNSYLVANIA

HEALTHCARE ADVOCATES, INC.,

      Plaintiff,

      v.                        Civil Action No. 05-03524
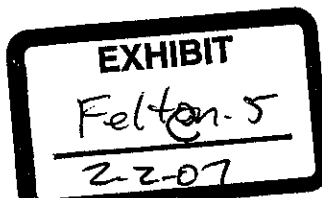
HARDING, EARLEY, FOLLMER & FRAILEY,

      Defendant.

## EXPERT REPORT OF EDWARD W. FELTEN

### I.    Introduction

1.    My name is Edward W. Felten. I am a Professor of Computer Science and Public Affairs at Princeton University, where I have taught for thirteen years. I also serve as the Director of Princeton's Center for Information Technology Policy.

2.    My research and teaching in computer science have focused on computer security and privacy, Internet software, computer systems, and technology law and policy. I have won awards for my work, including a National Young Investigator Award from the National Science Foundation and an Alfred P. Sloan Fellowship. In 2004, Scientific American magazine named me to the Scientific American Fifty, its

Expert Report of Edward W. Felten
Page 1

list of fifty global leaders in science and technology. A copy of my Curriculum

Vitae is attached to this report as Appendix A.

3.    I am being compensated at a rate of $700 per hour for my time working on

this matter, and my out-of-pocket expenses are being reimbursed.

4.    I have served as an advisor on information technology issues to many

agencies of the U.S. government, including the departments of Justice, Defense, and

Homeland Security; and the Federal Trade Commission. I have testified before

committee hearings of the U.S. Senate and U.S. House of Representatives.

5.    I was the lead technical expert witness for the Department of Justice in the

antitrust case *U.S. v. Microsoft*. In that capacity I testified twice, in December 1998

and June 1999, before Judge Thomas Penfield Jackson in the District of Columbia. I

testified briefly in 2000 in *Universal v. Reimerdes* before Judge Lewis Kaplan in the

Southern District of New York. I testified twice, in 2000 and 2003, in the patent

case *Eolas v. Microsoft*, before Judge James Zagel in the Northern District of

Illinois, Eastern Division. I testified in 2006 in *ACLU v. Gonzales*, before Senior

Judge Lowell Reed in the Eastern District of Pennsylvania.

6.    I have written and taught extensively about Internet software and Internet

protocols. For example, I have taught Princeton's course on computer networks,

which covers many Internet protocols.

7.    My research and teaching on computer security has included issues of

filtering, access control and circumvention. For example, I have taught about these

topics in courses on information security and information technology policy.

8.    I have been asked to render my expert opinion on certain technical issues relating to the claims in this matter, and to respond to the expert report of Gideon Lenkey.  My opinions are based on my training and experience, and on materials I have consulted.  I have consulted the following materials: all materials referenced in this report, and Web server log files from HCA and IA.

## II. Internet and Web Technology

9.    The Internet is a global network of computers constructed by patching together many local area networks that use widely varying communication media such as telephone lines, dedicated data cables, and wireless links.  The Internet is characterized by its global scope and by the use of certain standard data formats and protocols such as the Transmission Control Protocol and the Internet Protocol (together, "TCP/IP") that ensure that any two computers on the Internet can exchange information with each other.

10.    Computers on the Internet are named or addressed using two types of designators: IP addresses and DNS names.  IP[1] addresses are used by the Internet's infrastructure to route traffic to a computer, just as telephone numbers are used to route calls to a phone.   IP addresses are written in numeric form, such as 207.41.14.192.  DNS[2] names utilize a more user-friendly textual format; an example is www.uscourts.gov.  DNS names are translated automatically into IP addresses, so that users can take advantage of the friendlier DNS naming system while the

---

[1] IP stands for "Internet Protocol," which is the basic mechanism used to enable communication on the Internet.

[2] DNS stands for "Domain Name (or Naming) System," a system for giving computers on the Internet user-friendly textual names and translating those names into corresponding IP addresses.

Expert Report of Edward W. Felten
Page 3

Internet's basic infrastructure relies only on the more machine-friendly numeric IP addresses.

11.    The Internet serves as the basic communication infrastructure for a wide range of electronic information services and activities, including electronic mail, electronic discussion groups, teleconferencing, remote access by traveling workers to institutional databases, and the World Wide Web ("the Web"). Although the Web uses the Internet as its basic communication infrastructure, the terms "Internet" and "Web" are not entirely synonymous.

12.    The Web is characterized by a set of standard data formats, including HyperText Markup Language ("HTML"), and a set of standard communication protocols, such as HyperText Transfer Protocol ("HTTP"), that together allow computers to publish and view Web "pages" that may contain links to other such pages. The Web is made up of the global collection of pages that meet these format specifications, along with the global collection of computers that store and transmit pages on demand via these standard protocols.

13.    A Web browser is a computer program that allows a user, typically sitting at a personal computer or laptop, to access and read Web pages. Popular Web browsers include Microsoft's Internet Explorer, and the Mozilla project's Firefox.

14.    A Web server is a computer that stores Web pages and makes them available across the Internet using Web protocols such as HTTP. The term "Web server" is also used sometimes to refer to a software program that makes this possible.

15.    Web pages are identified by Web addresses such as

http://www.paed.uscourts.gov/contents.asp (the address of the table-of-contents page

for the Federal Courts in the Eastern District of Pennsylvania). These addresses are

technically known as Uniform Resource Locators ("URLs"). A URL designates a

Web server ("www.paed.uscourts.gov" in this case), a particular page on that Web

server ("contents.asp" in this case), and a protocol to use in retrieving the page

("http" in this case).

16.    Web servers normally keep log files, which record every request made to

the server. A Web server log consists of a sequence of entries. Each entry typically

records information about one request, such as the date and time, the specific page or

file requested, the address of the requesting computer, the outcome of the request

(success, or the type of error that occurred), and the amount of data transferred in

satisfying the request.

### A. Caching

17.    When a computer program accesses a Web page, the computer will

sometimes store a copy of that Web page, in case the page is needed again later.

This is known as "caching." If the program needs the page again, it can get it from

its local cache, rather than having to contact the Web server again to get another

copy of the page.

18.    Caching has advantages – for example, retrieving a page from the local

cache is faster than downloading it across the Internet – but it has drawbacks too.

One drawback is that if the page changes on the original server, the cached copy will

be out of date, and a program that relies on the cached copy will have outdated

information. For this reason, cached files are normally timestamped when they are initially put in the cache, and are discarded after some time interval has passed. For example, if cached files are discarded after twenty-four hours, then information retrieved from the cache will never be more than twenty-four hours out of date.

### B. Web Crawlers

19. A Web crawler is a computer program that catalogs Web pages. Crawlers try to discover as many pages as they can, and they download and store copies of the pages they discover. In other words, a crawler makes an archival copy of whatever portion of the Web it can discover.

20. Crawlers are used by search engines, such as Google and Yahoo, to help people find information on the Web. For example, if you type "Edward Felten" into Google, its response will point you to one of my Web pages. It can do this because it has archived the contents of that page (and many millions of other pages) and built an index of its archive, thereby allowing it quickly to find pages containing particular words.

21. Crawlers typically operate on their own, without a human sitting at the computer. The basic mechanisms of the Web, such as the HTTP protocol that is used for transferring most Web pages, are designed for computer-to-computer interaction. When a user is browsing the Web, his Web browser – a computer software program – is accessing Web pages on his behalf. In the same way, a program like a crawler can access Web pages, even if no human is present. (A human would have programmed the crawler (i.e., given it instructions) beforehand.)

Expert Report of Edward W. Felten
Page 6

22.   Crawlers are a well-known and accepted part of the Web "ecosystem." Many different crawlers exist.

## C. Robots.txt Files

23.   For various reasons, Web site publishers sometimes prefer that their sites, or certain parts of their sites, not be visited by crawlers.  Many of the people running crawlers are happy to comply with publishers' requests not to crawl their pages.

24.   In about 1993, technologists recognized that it would be useful to have a standardized method by which site publishers could ask crawlers not to visit their sites.  Ensuing discussions led to a consensus result known as the "robots.txt standard."  It is summarized in an online document entitled "A Standard for Robot Exclusion", available online at http://www.robotstxt.org/wc/norobots.html (hereinafter, "Standard for Robot Exclusion").

25.   This "standard" was not drafted or ratified by any formal standards body but was "ratified" only by common usage.  The Standard for Robot Exclusion summarizes its status as follows:

> It is not an official standard backed by a standards body, or owned by any commercial organisation. It is not enforced by anybody, and there no guarantee that all current and future robots will use it. Consider it a common facility the majority of robot authors offer the WWW community to protect WWW servers against unwanted accesses by their robots.

(Standard for Robot Exclusion)

26.   If a Web site publisher wishes to use the robots.txt mechanism, he creates a file called "robots.txt" and has his Web server offer this file for download.  This file, if present, contains a series of records. Each record gives the names of one or

more crawlers, and enumerates the parts of the web site that the specified crawlers are requested not to visit.

27.    The robots.txt file can be viewed by people, but it is intended to be read by crawlers. A crawler can look for the robots.txt file and (assuming it is present) can download and interpret it. The crawler can look for records that name it, and by reading those records it can learn that the web site author is requesting that it not visit certain parts of the site.

28.    If no robots.txt file is present, it is assumed that the Web site publisher welcomes crawlers.

29.    Many crawlers are programmed to read the robots.txt file and comply with the requests therein. Doing so is considered polite behavior.

30.    The state of a site's robots.txt file (i.e., whether the file exists, and if so what it contains) may change from time to time. The state at some point in time communicates the site publisher's requests as of that time. Requests are not retroactive. If I add a request to my robots.txt file at noon today, that request applies to crawler visits from noon today onward. It does not apply retroactively to requests made before noon today. It does not mean that crawlers' requests made before noon today were unwelcome, nor does it mean that crawlers should not store or use information gathered from my site before noon today.

31.    If a site publisher changes his mind about whether he wants a crawler to access his site, he can modify his robots.txt file accordingly. For example, if he is currently asking a crawler not to visit the site but he wants to welcome that crawler in the future, he can remove or modify the lines in the file regarding that crawler. A

crawler may check the robots.txt file from time to time to see whether the site publisher has changed his mind.

32.    A web site author can change his robots.txt file at any time, and the new version of the file is assumed to apply as soon as it is made available.  The contents of my robots.txt file as of noon today do not specify my wishes as of tomorrow morning -- to know my wishes tomorrow morning, you will have to look at my robots.txt file then.

33.    A Web site author who wants to exclude certain crawlers from accessing all or part or his site can do so by technical means.  For example, he can configure his Web server software to reject accesses (or accesses to certain parts of the site) that are labeled as coming from the targeted crawler(s)[3].

34.    The requests in a robots.txt file apply, by definition, only to the initial gathering or archiving of Web sites and pages by a crawler.  The requests do not apply to the subsequent storage, use, or redistribution of the archived sites and pages.

35.    The requests in a robots.txt file apply, by definition, only to crawlers. They do not apply to people.  Web site authors who want to exclude people from their sites use other mechanisms, such as password-protecting the site.

## III. Internet Archive's Technology

36.    Internet Archive ("IA") collects an archive of the Web, containing past versions of many Web sites.  To collect Web pages to store in its archive, IA uses crawlers.  IA's crawlers visit Web pages from time to time, recording their contents

---

[3] Requests for Web pages are labeled with a "User-Agent" header saying which computer program is making the request.  A Web server can reject requests if it sees the User-Agent header characteristic of a targeted crawler.

Expert Report of Edward W. Felten
Page 9

in IA's archive. IA's crawler is programmed to look for the robots.txt file on each site it visits, and to comply with any requests found therein.

37.    IA also provides an online service called the Wayback Machine[4] that lets users view archived versions of sites. A user can enter the URL (Web address) of a page or site into the Wayback Machine, and the Wayback Machine will then consult IA's archive to help the user view past versions of the requested page or site.

38.    IA has voluntarily chosen to have the Wayback Machine withhold archived copies of a site from users, if that site has a robots.txt file that asks the IA crawler not to visit that site. In doing this, IA goes beyond any request made by the site publisher in the robots.txt file – nothing in the robots.txt standard requires this from IA, nor does the standard indicate that doing so would be customary, polite, or otherwise desirable. A person familiar only with the robots.txt standard would have no particular reason to think that IA would do this, or that there would be anything unusual or improper about IA not doing it.

39.    To implement this behavior, the Wayback Machine would have to consult a site's robots.txt file every time a user tried to look at an archived version of that site. IA apparently intended the Wayback Machine to use caching, keeping copies of already-accessed robots.txt files for twenty-four hours. If implemented correctly, this would have caused the Wayback Machine to access sites' robots.txt files less frequently.

40.    During the relevant period in July 2003, the Wayback Machine's robots.txt caching mechanism was not working as intended. This is evident from the Web server logs of Healthcare Advocates ("HCA"), which show more accesses by

the Wayback Machine to HCA's robots.txt file than would have occurred had the caching mechanism been working as expected. Sometimes, when one would have expected the Wayback Machine to use a cached version of HCA's robots.txt file, the Wayback Machine instead retrieved the robots.txt file across the Internet from HCA's server.

41.    The number of extra accesses is only a few hundred, a very small number by Web standards. In any case, if HCA's robots.txt file were not changing over time, failure by the Wayback Machine to cache HCA's robots.txt file should not have affected the Wayback Machine's behavior (other than causing it to retrieve the robots.txt file more often). If HCA's robots.txt file were in place, properly constructed and unchanging, and if HCA's Web server were working correctly, then the Wayback Machine should have seen that same robots.txt file on every access. In other words, the Wayback Machine caching flaw described by Mr. Lenkey cannot by itself account for the Wayback Machine's delivery to Harding Earley of the archived files at issue here.

## IV. Harding Earley's Accesses to Internet Archive

42.    At issue in this litigation are a small number of accesses made by Harding Earley's employees to the Internet Archive web site.

43.    These were ordinary web page accesses, made via ordinary browsers. Harding Earley's employees did not use extraordinary, unusual, sophisticated, or devious methods to access these pages. They did not engage in "hacking". Mr. Lenkey does not point to any evidence that automated tools (beyond ordinary Web browsers) were used, nor am I aware of any such evidence.

---

[4] Though dubbed a "Machine", this service actually uses a group of computers.

Expert Report of Edward W. Felten
Page 11

44.     Even according to Mr. Lenkey, all that Harding Earley's employees did was to make ordinary accesses to publicly available Web pages. In other words, they simply asked their browser to fetch and display a web page –what every user does when browsing the web. In response to these requests, the Wayback Machine delivered the requested pages.

45.     Based on descriptions available on IA's web site, one might have expected the Wayback Machine to refrain from delivering the requested pages. Nevertheless, the Wayback Machine did deliver the pages.

46.     I agree with Mr. Lenkey that the reasons the Wayback Machine provided the pages cannot definitively be determined now (Lenkey Report at p. 6: "for reasons unknown").

47.     We can, however, identify several likely scenarios that could have led to the observed results. First, HCA might have tinkered with its robots.txt file, so that it either did not deny access to the IA crawler or was improperly formatted in a way that confused the Wayback Machine servers. Second, HCA's web server might have been misbehaving, failing to produce the correct robots.txt file in response to the Wayback Machine's requests for it. Third, the Wayback Machine might have had trouble reading or interpreting the contents of the robots.txt file.

48.     These three examples do not exhaust the possibilities, but they do give an idea of how errors by HCA or by IA could have led to the observed behavior. Regardless of what happened, Harding Earley's role was simple – to request files by the usual means, and to receive them.

49.    Mr. Lenkey claims that Harding Earley's accesses to the Wayback Machine were unauthorized and that the unauthorized accesses were intentional. I disagree with his characterization of these accesses.

50.    His characterization appears to rely on an assumption that Harding Earley's employees knew, or could reasonably be expected to have known, that the Wayback Machine was supposed to be preventing them from accessing the files that it was allowing them to access. This assumption seems to rest on two assertions, regarding Harding Earley seeing HCA's robots.txt file, and Harding Earley seeing an "query exclusion page" supplied by the Wayback Machine. According to Mr. Lenkey, a person seeing these two files should have known that future accesses to the requested files would be unauthorized.

51.    Regarding the robots.txt file, HCA's server logs indicate that somebody at Harding Earley saw some version of HCA's robots.txt file during the relevant period, but the record does not tell us what exactly was in the version of the robots.txt file that Harding Earley saw. In any case, as explained above in paragraph 35, the robots.txt file at most requested that *IA's robot* not visit HCA's site. HCA does not assert that the robots.txt file requested that *Harding Earley* refrain from accessing HCA's site. In any case, the robots.txt file, whatever it contained, could not have denied permission to any *person*.

52.    According to Mr. Lenkey, the relevant portion of HCA's robots.txt file consisted of these two lines of cryptic text:

User-agent: ia_archiver
Disallow: /